

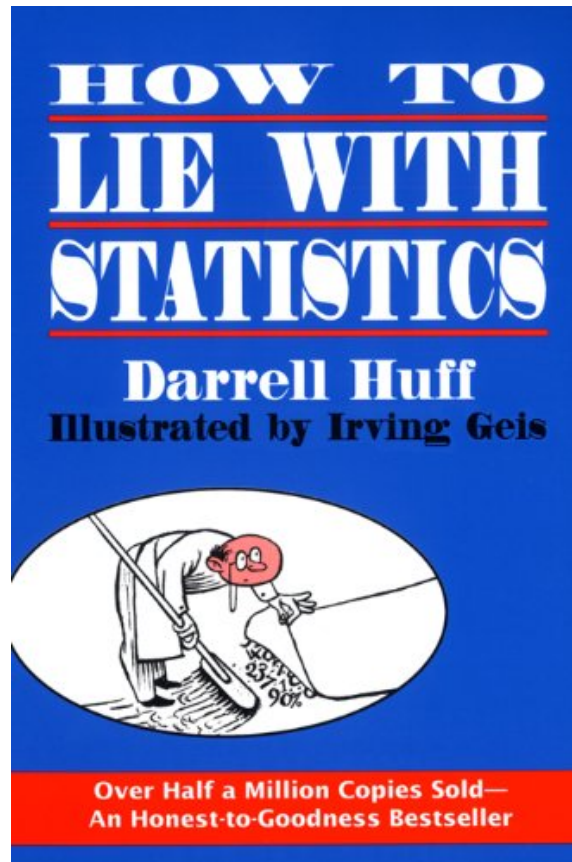
Statistics (I)

SLIDES BY:

JIANNAN WANG

<https://www.cs.sfu.ca/~jnwang/>

Why Should You Care?



“ There are three kinds of lies:
lies, damned lies, and statistics ”

1. <u>The Sample with the Built-in Bias</u>	13
2. <u>The Well-Chosen Average</u>	29
3. <u>The Little Figures That Are Not There</u>	39
4. <u>Much Ado about Practically Nothing</u>	55
5. <u>The Gee-Whiz Graph</u>	62
6. <u>The One-Dimensional Picture</u>	68
7. <u>The Semiattached Figure</u>	76
8. <u>Post Hoc Rides Again</u>	89
9. <u>How to Statisticulate</u>	102
10. <u>How to Talk Back to a Statistic</u>	124

Simpson's paradox

Is UC Berkeley gender biased?

	Applicants	Admitted
Men	8442	44%
Women	4321	35%

~~YES!~~

Simpson's paradox

Is UC Berkeley gender biased?

Department	Men		Women	
	Applicants	Admitted	Applicants	Admitted
A	825	62%	108	82%
B	560	63%	25	68%
C	325	37%	593	34%
D	417	33%	375	35%
E	191	28%	393	24%
F	373	6%	341	7%

NO!

Women tended to apply to competitive departments with low rates of admission

Outline

Statistical Thinking

Descriptive Statistics

Inferential Statistics

Outline

Statistical Thinking

Descriptive Statistics

Inferential Statistics

Statistical Thinking

1. Data is just a **sample**
2. Your goal is to infer a **population**
3. Think about how to go “backwards” from the **sample** to the **population**

Example 1. Image Classification

Is it a dog or a cat?



Dataset: 1000 images collected from the Web

Without Statistical Thinking

Treat the 1000 images as the population

- > Train a model on the data
- > Evaluate a model on the same data
- > Model accuracy: 95%

With Statistical Thinking

What is the population?

- All the images in the Web

What is your dataset?

- A sample of 1000 images drawn from the Web

What should you do?

- Split the dataset into a training dataset and a test dataset
- Train the model on the training dataset
- Evaluate the model on the test dataset

Example 2. Poll Prediction

Who will win the election?



Dataset: A survey of 100 people

Without Statistical Thinking

Treat the 100 people as the population

- > Count the number of people who wants to vote for Hillary, e.g., 52
- > Count the number of people who wants to vote for Trump, e.g., 48
- > Hillary will win the election

With Statistical Thinking

What is the population?

- All the people who will vote in the election day

What is your dataset?

- A sample of 1000 people before the election day

Analysis result

Hillary: 52% \pm 3%
Trump: 48% \pm 2%

Assumption: People have not changed their votes since the time of the poll

Summary

Statistical Thinking

- Sample, Population and Their Connection
- With vs. Without Statistical Thinking

Descriptive Statistics

Inferential Statistics

Outline

Statistical Thinking

Descriptive Statistics

Inferential Statistics

Descriptive vs. Inferential Statistics

Descriptive Statistics: e.g., Median

- Why? Aim to understand the data
- How? Data summarization, data visualization, etc.

Inferential Statistics: e.g., A/B Testing

- Why? Aim to use the data (i.e., sample) to learn about a population
- How? Estimation, confidence intervals, hypotheses testing, etc.

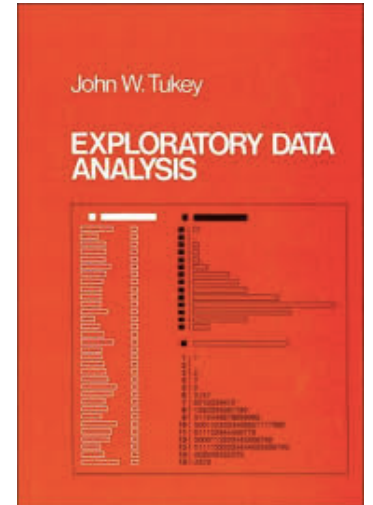
Exploratory Data Analysis (EDA)

The process of doing descriptive statistics



John W. Turkey

- Professor at Princeton University
- Founding chairman of the Princeton statistics department in 1965
- Worked on EDA at Bell Labs since 60's
- Wrote a book entitled "Exploratory Data Analysis" in 1977



EDA is like detective work



From John Turkey

“ Exploratory data analysis is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those that we believe to be there ”

EDA with Dataprep

EDA Solutions in Python

Solutions	APIs	Efficiency	Out-of-memory Data
Pandas + Matplotlib ¹	Designed for plotting	Slow	Not supported
Pandas-profiling ²	Designed for profiling	Slow	Not supported
Dataprep.eda ³	Designed for EDA	Fast	Supported

1. <https://pandas.pydata.org/pandas-docs/version/0.13/visualization.html>
2. <https://github.com/pandas-profiling/pandas-profiling>
3. <http://dataprep.ai/>

EDA Tasks vs Dataprep API

Task 1. Understand the distributions of all columns

- `plot(df)`

Task 2. Dive into the distribution of a single column

- `plot(df, x)`

Task 3. Dive into the relationship between two columns

- `plot(df, x, y)`

EDA Tasks vs Dataprep API

Task 4. Understand the missing values of all columns

- `plot_missing(df)`

Task 5. Dive into the missing values w.r.t. a single column

- `plot_missing(df, x)`

Task 6. Dive into the missing values between two columns

- `plot_missing(df, x, y)`

EDA Tasks vs Dataprep API

Task 7. Understand the correlation for all pairs of columns

- `plot_correlation(df)`

Task 8. Dive into the correlation w.r.t. a single column

- `plot_correlation(df, x)`

Task 9. Dive into the correlation between two columns

- `plot_correlation(df, x, y)`

Correlation Analysis

Correlation

- It is a measure of relationship between two variables

Why is correlation analysis useful?

- For understanding data better
- For making predictions better

Case Study: How to do correlation analysis

Height and weight are correlated

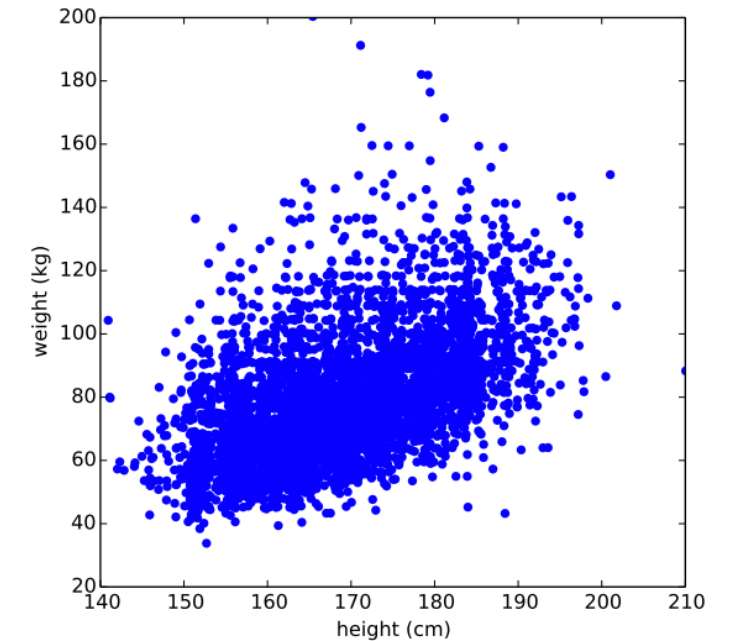
1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0

Source: *Think Stats -- Exploratory Data Analysis in Python*

Idea 1. Visualization

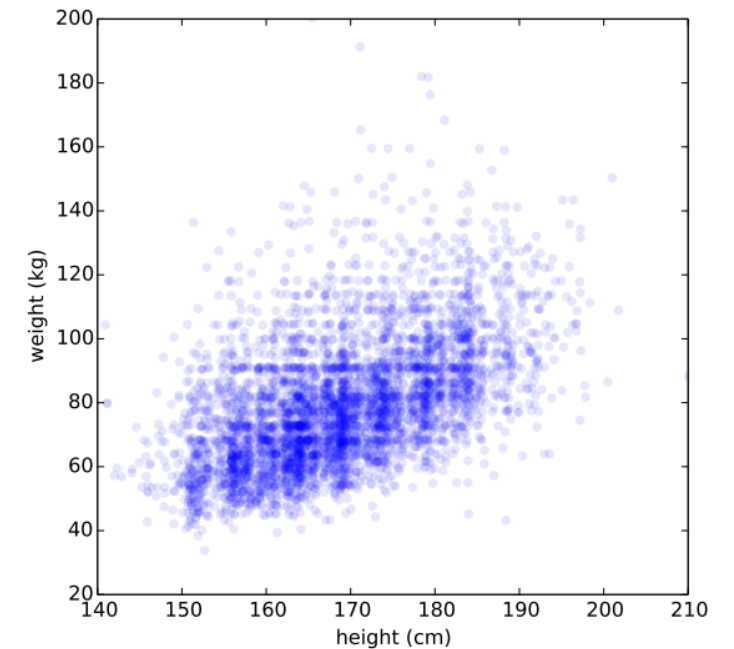
Scatter Plot

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



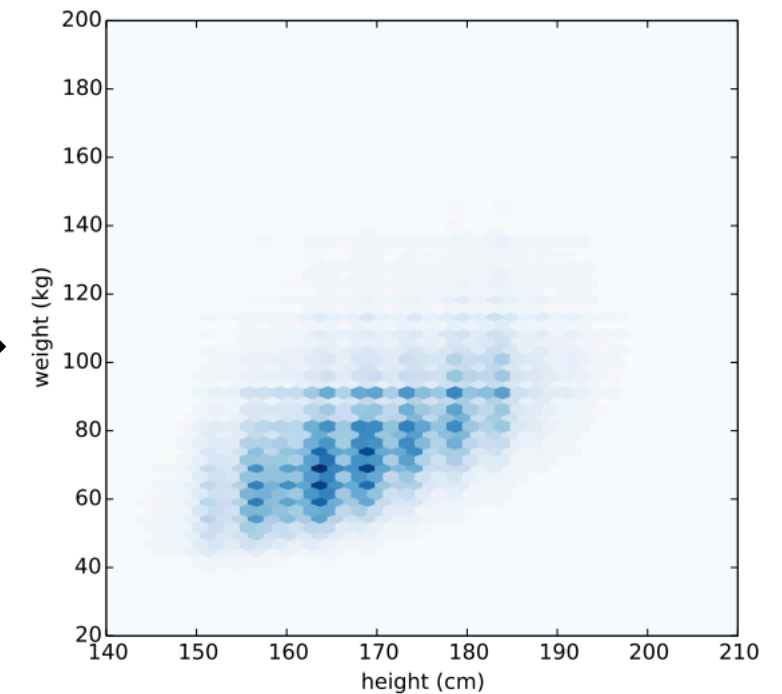
Scatter Plot (with transparency)

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



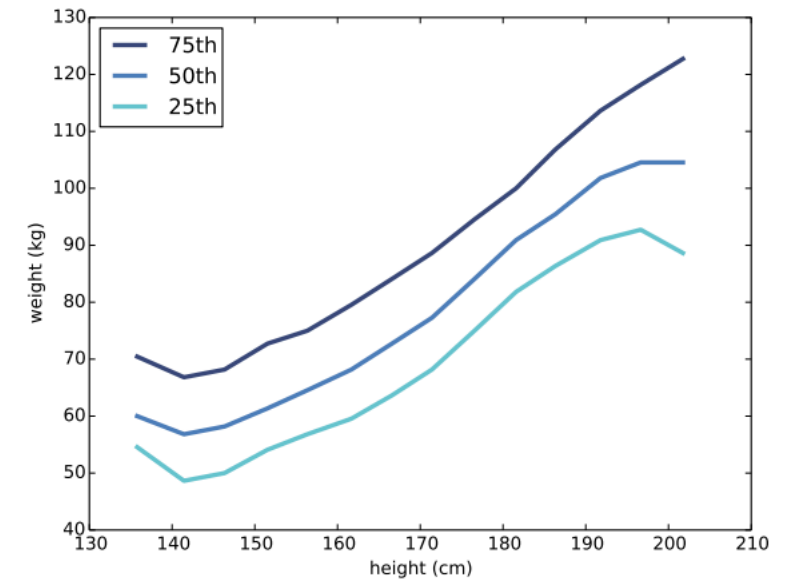
Hexbin Plot

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



Characterizing relationships

1	height	weight	age	male
2	151.765	47.8256065	63	1
3	139.7	36.4858065	63	0
4	136.525	31.864838	65	0
5	156.845	53.0419145	41	1
6	145.415	41.276872	51	0
7	163.83	62.992589	35	1
8	149.225	38.2434755	32	0



Idea 2. Correlation Coefficient

Covariance

Covariance is a measure of the **tendency** of two variables to vary together.

$$\text{cov}(X, Y) = \text{E} [(X - \text{E}[X])(Y - \text{E}[Y])]$$

$$\text{cov}(X, Y) = \text{E}[XY] - \text{E}[X] \text{E}[Y]$$

Hard to interpret
113 kilogram-centimeters

Pearson's correlation

Pearson's correlation is a measure of the linear relationship between two variables

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

Easy to Interpret

- $[-1, 0) \rightarrow$ Negative Correlated
- $(0, +1] \rightarrow$ Positive Correlated
- -1 or $+1 \rightarrow$ Perfectly Correlated

What about non-linear relationship?



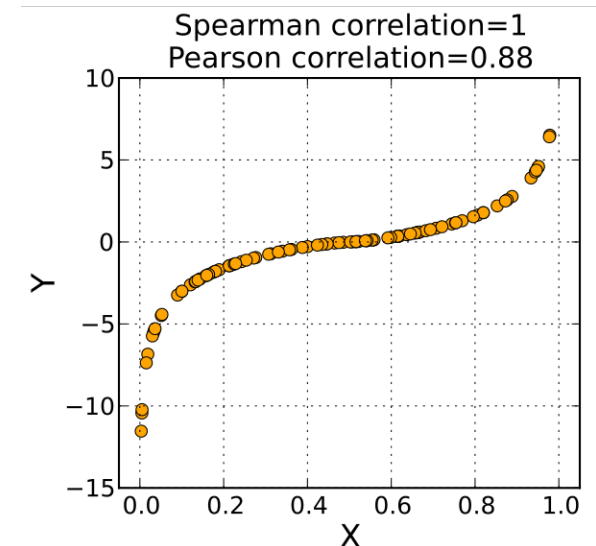
Spearman's rank correlation

Spearman's rank correlation is a measure of monotonic relationship between two variables

$$r_s = \rho_{r_X, r_Y} = \frac{\text{cov}(r_X, r_Y)}{\sigma_{r_X} \sigma_{r_Y}}$$

Advantages

- Mitigate the effect of outliers
- Mitigate the effect of skewed distributions



Summary

Statistical Thinking

Descriptive Statistics

- Descriptive vs. Inferential Statistics
- Exploratory Data Analysis with Dataprep
- Correlation Analysis

Inferential Statistics

Outline

Statistical Thinking

Descriptive Statistics

Inferential Statistics

- Estimation

Estimation

Problem statement

- Estimate a numerical value associated with a population

Examples

- Estimate the percentage of the people in the US who will vote for Trump
- Estimate the median annual income of all households in the US

Example: Median Annual Income

How to estimate the median annual income of all households in the US?

- Randomly select 10,000 households from the US
- Report their median annual income: 50,000USD
- BUT, we need to report something like

50,000 \pm 500 USD

A Naïve Solution

- Randomly select 10,000 households from the US
- Report their median annual income

Repeat this process for
100 times

50,000 49,600 50,200 ... 49,200

You have to survey 1,000,000 million households in total 😞

A Smart Solution: Bootstrapping

Key Idea: Resampling

- Sample with replacement from the original data sample

Population: 1, 1, 8, 2, ... 3, 3

Sample: 3, 8, 1, 8, 3

Resample: 8, 3, 3, 3, 1

A Smart Solution: Bootstrapping

- Randomly select 10,000 households from the US
- Draw a resample from the 10,000 households
- Report the median annual income of the resample

Repeat this process for
100 times

You do NOT need to survey any new household. 😊

Notes on Bootstrapping

Start with a large random sample (at least 30)

Replicate the resampling procedure as many times as possible (more than 1000 times)

Does not work for min/max

Conclusion

Statistical Thinking

- Sample, Population and Their Connection
- With vs. Without Statistical Thinking

Descriptive Statistics

- Descriptive vs. Inferential Statistics
- EDA with Dataprep.eda
- Correlation Analysis

Inferential Statistics

- Estimation and Bootstrapping